

Open Refine

Open Refine, formerly known as Google-Refine, is an open source data cleaning & data transformation tool that provides a user-friendly interface for cleaning, transforming & enhancing messy or unstructured data.

Useful to identify similar redundant rows & merge them into clusters.

Image Labeling . PIL

Python Imaging Library, is a popular library for basic image processing tasks, but it has been succeeded by the Pillow library, which is a more actively maintained and enhanced fork of PIL.

Image opening &
saving.

Image manipulation

Color &
channel manipulation

Basic
Drawing.

Wadker

Week 4

MODEL THE DATA

which software to use when & for which purpose!

- Excel . for i) Correlation.

Date:

Page No.:

(How often is there an effect)

(Eg: How often increase in Math marks imply increase in Phy marks?)

- ii) Regression.

(How much is there an effect?)

- iii) Outlier Detection.

(When does this relation fail?)

- Python for i) Classification

(Given data pt, which grp should it belong based on the history data?)

- ii) Forecasting.

(Future prediction)

- iii) Clustering

(can we group things into similar entities?)

- Others

► R/R studio one of the most powerful tools for data modeling

► Rattle non-programmatic R-based application that help you analyze data

► Py caret similar to Python; helps you build models and explore models fairly, quickly.

* Correlation (using Excel). Use Data > Data Analysis (Add if not there using Add-in)

- Interpret correlation matrix.

	A	B	C
A	1		
B	ab	1	
C	ac	bc	1

$$1 > (ab) > 0$$

Correlation

Date:	
-------	--

+ve correlation

slope { $ab \rightarrow 1$ } high/strong +ve corr
of linear map { $ab \rightarrow 0$ } weak +ve corr.

$$-1 < (ab) < 0 \quad -ve \text{ correlation}$$

$(ab) \rightarrow -1$ strong -ve corr

$(ab) \rightarrow 0$ weak -ve corr.

- Use scatterplots to visualize each correlation

* Regression (using Excel).

Use Data > Data (Add if not present)
Tab Analysis > Regression.

- Single Linear Regression.

1 indep var.

1 dep var.

single linear regression
(automatic selection based on selected data)

- Multiple Linear Reg.

(> 1) indep; 1 dep.

containing 1 or ≥ 1 indep var).

SUMMARY OUTPUT

Mult R Square 0.884585151

R square 0.7824...

Adj R Sq 0.7740...

for single linear regression

for multiple linear regression

Coef.

Intercept c_1

x_1 City Pop. c_2

x_2

$$y = c_1 + c_2 x_1 + c_3 x_2 + \dots + c_n x_n$$

* Outlier Detection .

	A	B	
1	Values	Outliers	
2	33090	= FALSE	
3	23671	FALSE	
4	255215	TRUE	
	:	:	

$$= \text{OR}(A2 < D1, A2 > D5)$$

Date :

Page No.

	C	D .
1	Q1	$3453 = \text{QUARTILE.EXC}(A:A, 1)$
2	Q3	$100158 = \text{QUARTILE.EXC}(A:A, 3)$
3	IQR	$96705 = D2 - D1$
4	Lower B - (...) → 0 (can't be zero)	
5	Upper B 245215.5	

$$= D1 - (1.5 * IQR)$$

$$(Q1 - 1.5 * IQR, Q3 + 1.5 * IQR) .$$

$$\rightarrow = D2 + (1.5 * IQR)$$

Model the Data : Pycaret library, Shap library.

Pycaret is an open source Python library designed to simplify the end-to-end machine learning workflow. It was created to help data scientists & machine learning practitioners streamline the process of building and deploying machine learning models by automating many of the common tasks & reducing the need of extensive coding.

```
import pycaret
```

```
from pycaret.datasets import get_data
```

```
index = get_data('index')
```

displays the 30 different datasets present in

pycaret-dataset

juice data.shape
returns the dimensions of the dataset from pycaret.classification import *
clf1 = setup(data=juice, target='Purchase',
log_experiment=True, experiment_name='juice',
normalize=True, feature_selection=True)

random seed. to ensure some randomization

we setup an experiment for pycaret to use
this data and get a sense on what kind
of model are we dealing with.

The output you get; based on this you can have
quick analysis of stuff for further process
Plug & play PyCaret.

training_data = get_config(variable="X-train")

we receive the training dataset

training_data.shape

models()

shows the details of models performed on the
dataset, and how well is it to be chosen.

best_model = compare_models()

best_model

shows the best-fit model.

Either use this or the user may explicitly

use some other model of their choice

say lr = create_model('lr') ~~for now~~

tuned-best-model = tune-model(best-model)

Fine tunes the hyperparameters of the model to
establish better performance

Date :

Page No.

plot-model(tuned-best-model)

displays ROC curve

Receiver Operating Characteristic

plot-model(tuned-best-model), ~~plot='confusion matrix'~~

plot-model(tuned-best-model, plot='feature')

different plot curves.

interpret-model(tuned-best-model)

shows the graph with shaft values.

save-model(tuned-best-model, ~~model-name='best-model'~~)

saves the model as a pickle file

extension = .pkl.

* Pycaret helps a lot in automating Machine Learning.

Clustering with Python.

import pandas as pd

import numpy as np

clustering method KMeans gets imported
from sklearn.cluster import KMeans

stockData = pd.read_csv("https://raw.githubusercontent.com/.../
...csv")

StockData.head()

Multiplying those % columns that range from 0 to 1
stockData['Free Float...'] = stockData['Free ...'] * 100
stockData['ROE %'] = stockData['ROE %'] * 100
:
StockData.describe()

Date :

Page No.

features = stockData[1], columns [2:]
storing features.

from sklearn.preprocessing import MinMaxScaler
scaler = MinMaxScaler()

stockDataFeatures_scaled = scaler.fit_transform(stockData[features])

stockDataFeatures_sorted = pd.DataFrame(stockDataFeatures_scaled,
 columns=features)

Scaling the DB

stockDataFeatures_scaled.describe()

stockDataFeatures_scaled

kmeans = KMeans(7) # stores the model object

clus = kmeans.fit_predict(stockDataFeatures_scaled)

stockData['cluster'] = clus

stockData['cluster'].value_counts()

clusterDesc = pd.DataFrame(stockData.iloc[:, 2].groupby('cluster').mean().round(3))

clusterDesc.insert(0, 'size', stockData['cluster'].value_counts())

•) What is a shapefile and how to read a shapefile using python?

→ extension (.shp)

A shapefile is a widely used geospatial vector data ~~as~~ format for Geographic Information Software (GIS).

A shapefile stores geometric & attribute data for geographic features, such as points, lines & polygons & is often used to represent geographic info like maps, boundaries & locations.

import geopandas as gpd

shp-path = 'pathTo/your-shapefile.shp'

gdf = gpd.read_file(shp-path)

printing head...

gdf.head()

Geospatial Analysis. also known as spatial analysis or geographic information system analysis (GIS), is the process of examining & interpreting data that has a geographic or spatial component.

W5
Deep learning models

code.

For this example, we are checking ~~there~~ whether # there is a need of a new store near The Empire

State Building

Geospatial analysis :

The objective of using geospatial analysis is to make decisions about opening up new restaurants (or retail stores, bank branches, airports etc). Here in this example, we will use a dataset consisting of Starbucks & McDonald's locations in New York.]

import pandas as pd

import numpy as np

import folium

import geopy.distance

Read file containing location details

df = pd.read_csv('https://drive.google.com/---')

df.head(10).

	lat	long	store	address
0	40.7	-74.006	Starbucks	375 Pearl St, New York, NY 10038
1				

Merge lat & long columns to help us
df['coordinate'] = '(' + df['lat'].astype(str)
+ ',' + df['long'].~~date~~.^{str}.~~astype(str)~~ + ')'

Page No.

df.head(10).

Search from web : The lat-long of the target

location

Empire State Building : (40.7188, -73.9852)

NY_coord = (40.748488, -73.985238)

Step 1: Compute distances from the stores

~~distance []~~

distances_km = []

for row in df.iterrows(index=False):

distances_km.append(geopy.distance.distance
library method)

(NY_coord,
row.coordinate).km)

df['Distance'] = distances_km

df.head(10)

Visualize Data on map (Step 2)

Empire State building coordinates

m = folium.Map(location=[40.748488, -73.985238],
zoom_start=10)

Place markers for stores in the map. → returns (ind, row)

for i, row in df.iterrows():

lat = df.at [i, 'lat']

lon = df.at [i, 'lon']

store = df.at [i, 'store']

Date :

Page No.

If store == 'McDonalds' :

color = 'red'

else:

color = 'green'.

folium.Marker (location = [lat, lon], popup = store,
icon = folium.Icon (color = color)).

add_to(m)

m. # displays the map

Step 3 # Compute no of stores located in a given radius

All stores in a radius of x=10 kms

df.loc [df ['Distance'] > 10]

'loc' accessor
that allows
to select data
from a DataFrame by
label-based indexing.

Step 4 : Compute the closest & farthest store from

target_loc.

```
df-farthest = df.iloc [df.groupby('store')  
                      .pd.Series.idxmax]  
df-closest = df.iloc [df.groupby('store')  
                      .pd.Series.idxmin]  
df-new = pd.concat ([df-farthest, df-closest])  
df-new # prints the dataframe with closest & farthest
```

Visualize the farthest & closest using folium

similar to what has been before

Sentiment Analysis with Excel & Azure ML

* Azure ML, Azure ML

Open Excel.

Insert > My-add-ins > Azure ML > ^{Text sentiment} analysis

Then select data

Select output data

* TextBlob library

```
import pandas as pd
```

```
import numpy as np
```

```
import spacy
```

```
from textblob import TextBlob
```

```
from sklearn.metrics import classification_report
```

importing all necessary libraries.

```
data = pd.read_csv('...csv')
data.head() # reading csv files & showing head data.head()

data['TextBlob-Subjectivity'] = data['review'].apply(lambda x:
# applying TextBlob subjectivity formula TextBlob(x).sentiment.
subjectivity)

data['TextBlob-Polarity'] = data['review'].apply(
# applying TextBlob Polarity (lambda x: TextBlob(x).sentiment.
polarity)
formula)

data['TextBlob-Analysis'] = data['TextBlob-Polarity'].apply(
# applying TextBlob-Analysis formula. (lambda x: 'negative' if x < 0
else 'positive')
```

data

```
print(classification_report(data['sentiment'],
data['TextBlob-Analysis']))
```

Image classification.
with Keras, with Google Cloud Platform.

- Keras: Keras is an open-source deep learning framework written in Python. It is designed to be user-friendly, modular and extensible, and making it accessible for beginners while providing advanced functionality for experienced practitioners.

```
from keras.models import Sequential  
from keras.layers import Dense
```

Date:

Page No.

Define sequential model

```
model = Sequential()
```

add layers to the model

```
model.add(Dense(units=64, activation='relu', input_dim=100))
```

```
model.add(Dense(units=10, activation='softmax'))
```

compile model

```
model.compile(loss='categorical_crossentropy', optimizer  
= 'sgd', metrics=['accuracy'])
```

Train the model

```
model.fit(x_train, y_train, epochs=10, batch_size=32)
```

W6

Design
your
output

Visualizing forecasts with Excel.

Given a data series (time-series) of currency
value fluctuations

- Trendlines? 'Insert' → 'Line Chart'.
(sparkline)

• correlation visualization := Scatter Plots.
(Add trendline to understand the linear regression relation)

Date:

Page No.

• Correlation Matrix 'Data Analysis' → 'Correlation'.

• colour grading based on values 'Conditional formatting' → 'Color Scales'

Visualization of animated data with flourish.

Use 'Bar Chart Race'

Kumu.

(Primarily) to visualize ~~the~~ interlinkages.

105
Design
your
output

Design the output

- General
1. Excel
 2. Google Data Studio
 3. Power BI
 4. Tableau

Specialized

1. Excel (VBA)
2. Flourish studio
3. Kumu
4. QGIS.

useful for
network
diagram

+ animated
includes visualizations

Geospatial analysis

Visualization with Excel

- ✓ Pivot Table
- ✓ Trends using Sparklines
- ✓ Correlation matrix.
- ✓ Charts from pivot table.

Visualization with Powerpoint.

- ✓ Use pivot table and then scale the unit such that it can taken as the width of a box in ppt.

- ✓ Insert box & textbox.
- ✓ Have a slide for each year/category
- ✓ Use transition 'Morph' & do some editing.

Visualization with Kumu

Kumu is a web-based platform designed for visualizing & exploring complex systems, & relationships.

Page No.

Key features : Relationship Mapping

Data Import & Integration

Visual customization

Dynamic Mapping .

Flourish

Visualization with ~~Flourish~~ Studio .

Flourish Studio refers to a web-based data visualization platform called Flourish .

Flourish is designed to make it easy for users to create interactive and engaging data visualizations without requiring extensive coding or design skills .

✓ user friendly interface

✓ Templates & designs

✓ Data import & connectivity

✓ Interactivity .

(W7)

Narrate a story

Numbers

suppose usage of some numerical statistics, percentages, etc

Date:

Page No.

Visuals

Charts, trendlines, etc

Animated charts.

Text

Descriptive statistics

Illustrations

Pictographical representation.
(comics)

- narrativeschool.com/quill/tableau/free-trial

(W8)

Deploy the results

Anonymize

the data

Build

app

Host

app

- Tools (ARX, Amnesia)
- Libraries (Faker, namesis)
- Notebook (colab, kaggle)
- Data app (streamlit, shiny)
- Web app (Flask, Tomado)
- Content (GitHub, Desktop)
- Apps (Heroku, Glitch)
- Infra (AWS, Azure)

Secure App → Scale App

libraries to build web applications

- Streamlit

Streamlit is an open source Python library that is used for creating web applications for data science & machine learning. It is designed to be easy to use, and allows developers to create interactive and customizable web apps with minimal effort.

Basic Streamlit components

- 1) Text
- 2) Interactive Widgets
- 3) Charts & plots (matplotlib.pyplot)
- 4) Dataframe display (pandas)
- 5) Interactive maps (folium lib)

ngrok

ngrok is a tool that creates secure tunnels to localhost, allowing you to expose a local server to the internet.

Services to host web application.

Heroku Heroku is a cloud platform as a service (PaaS) that allows developers to deploy, manage and scale applications easily. It provides a platform for building, deploying, and scaling web applications without the complexity of managing the underlying infrastructure.

22nd Dec

- TDS by Q
- MLP by Q
- MLP session

23rd Dec

- MLP, TDS

Date : revision
Page No.

22/12/23

PYQ

TDS Apr 2023

Q1) The dataset consists of geographic, demographic, information about countries & their respective GDPs. You would like to visualize this data & study the relationship between the location of countries & their GDPs. You decide to use PowerBI to visualize the dataset. But you would also like to generate summary of the data.

Suggest an appropriate way.

1) Map visualization

2) Scatter plots

3) Treemaps

4) PowerBI Q&A

5) Card Visuals

6) Key Influencers

Q2) Your project requires you to study the districts & their respective health indicators. You have a shapefile with you ~~that~~ provides the required details. The objective of the project is to identify and carve out ~~districts~~ districts that present high levels of - health indicators. How could QGIS be used ~~to~~ here?

QGIS: Quantum Geographic Information System.

It is a software, ~~an~~ open source software, that allows users to view, edit & analyze geospatial data.

Page No.

— Here QGIS can be used to create shapefiles for districts with high levels of health indicators.

Q3) What are the two outputs provided by Excel Azure Machine learning plugin?

Sentiment, Score.

Q4) The dataset consists of year, annual cotton production, annual rainfall, loan interest rates & fuel prices. You would like to compute the correlation coefficient between annual cotton production & other variables in the dataset to analyse the effects of various variables on the target var.

Using Excel, it is a wonderful way to ~~get~~ get correlation coefficient.

Q5) You would like to prepare your dataset before analysis. You choose pandas-profiling library to perform exploratory analysis.

Choice is appropriate

Pandas profiling provides information about outliers

Q6) Comicgen is a useful tool in narrating Data stories using comics. Which of the following is ~~not~~ a function of comicgen?

- ✓ Comicgen creates comic characters
- ✓ Comicgen provides options to custom create different comic characters & their emotions and pose.
- ✓ Comicgen can be easily integrated into Google sheets or Excel to narrate your data stories.

Q7) A very large Matrix A has a lot of zeros in it. Which function from scipy library is useful for efficient storage of A?

CSR-matrix → compressed sparse row matrix.

Q8) Which library has functions & tools that are useful in analysis of large graphs?

scikit-network

Q9) Kunuu is a tool that allows you to:

visualize complex network data.

Q10) Which library is used to extract data from wikipedia pages? wikipedia

Date :

Page No.

Q11) A dataset is provided to you about countries & respective populations. You plan to visualize the data in Tableau using map representation. But you are unable to do so because the map representation is not activated for you to choose. What could be the issue?

col-Name col-type

Country String

Population Integer

There might be column type incompatibility issues

Q12) _____ is helpful to understand the structure of a website before writing a scraping script.

Developer Tools

Q13) requests library has to get a webpage's html contents into Python.

Q14) type of location can be retrieved using Nominatim in Python. TRUE

Q15) Is there any restriction to the type of delimiter that can be used in text-to-column func' in Excel? No, no restriction

Date:

Page No.

Q16) What is the y-axis in autocorrelation plot?
correlation
x-axis in autocorrelation plot is time lag.

Q17) Which tools cannot be used for anonymising the data? PowerBI

Q18) For a one-time anonymization, static anonymization is sufficient. TRUE

Q19) Variable X has values AA, BB, CC. This info is represented as shown below:

AA BB CC

0 0 1

1 0 0

pandas.get_dummies

0 1 0

1 0 0

Q20) k-means is typically influenced by the start values. What option in sklearn.cluster.KMeans helps reduce the impact?

n_init

Q21) Pandas dataframe has salary-range
 You are interested in
 finding out how many
 employees are present in each
 category?

high	
medium	
medium	
low	
high	

DF['salary-range'].value_counts()

dataset Q22) We have an imbalanced class. Which feature in DecisionTreeClassifier() will help us tackle the problem?
 class_weight.

Q23) $\sum_{i=1}^{100} |y_i - \hat{y}_i| = ?$ sklearn.metrics.
 mean_absolute_error

Q24) We are interested in fitting an ARIMA model to our time-series data. Specifically we are interested in a moving average model of 0, setting a lag value of 4 for autoregression, & a difference order of 1.

ARIMA(..., order(4, 1, 0))
 ↓ ↓ →
 AutoRegressive lag value moving
 Integrated differorder avg model
 Average

Q25) pycaret is a ?
pycaret is a low-code machine learning library.

Date :

Page No.:

Q26) 'subjectivity' & 'polarity' are two properties returned by sentiment function of which library?
TextBlob .

Q27) what is the purpose of 'subjectivity' & 'polarity' ?

'subjectivity' score ranges from 0.0 to 1.0

very objective

'polarity' score

ranges from -1.0 to 1.0

very subjective

negative
sentiment

0.0
positive
sentiment

Q29) What does classification_report function from the sklearn.metrics module?

It builds a text report displaying the main classification metrics

Page No.

Q30) csr-matrix from scipy library

helps reduce matrix space when there are a lot of zero entries in the matrix

Q31) Google Studio is a tool that allows you to create dashboards for small scale projects

Q32) what is the name of the tab that is used to identify API calls in the Inspect element in any browser?

Network

Q33) Which of the following libraries are used to construct API urls? urllib

Q34) If a time series with a significant autocorrelation at lag 1 implies that

current observation is correlated with previous observation.

Q35) Sentiment function in TextBlob calculates a numerical score for the sentiment of the text.

Q36)

Match the following

1. Worksheet

a) contains a seqⁿ of worksheets or dashboards that work together to convey information
Page No.

2. dashboard

b) contains single view along with shelves, legend, data pane.

3. story

c) collection of view from multiple worksheets

1 → b

2 → c

3 → a

Identify the join used below:



Inner join.

Q38)

Which attribute does not belong to the category of data while categorizing data based on specifics of the disclosure risks from which a dataset is to be ~~protected~~ protected?

non-identifying attributes.

Q39)

What is one of the ways / methods of removing outliers?

IQR analysis

Q40) On applying OHE, what happens to the dataframe
the #columns increases.

Q41) _____ table in Excel is a way of quickly
summarizing your data. Pivot table

Q42) Fill the missing line.

distances_km = []

for row in df.iterrows(index=False):

distances_km.append(

... 'missing line')

)
df['distance'] = distances_km.

df.head(10)

geofy-distance.distance(location, location2).km

Q43) Write a piece of code to display details of
9 scheduled airlines using BeautifulSoup.

import requests

import pandas as pd

from bs4 import BeautifulSoup

website_url = requests.get('...').text

soup = BeautifulSoup(website_url, 'html.parser')

required_table = soup.find_all('table')[0]

```
df = pd.read_html(str(required_table))
df = pd.DataFrame(df[0])
df
```

Date : _____
Page No. _____

Q44) Write the correct format of logical expression in Tableau.

IF [Sales] > 1000 THEN 'High Sales'
ELSE 'Low Sales' END

Q45) what is the attribute to get latitude & longitude of a location

Normalise().geocode('location-name').latitude
longitude

Q46) In streamlit, how to create a dropdown menu select for hobbies having 'Dancing', 'Reading', 'Sports' as options.

```
import streamlit as st
st.selectbox("Hobbies:", ['Dancing', 'Reading', 'Sports'])
```

Q48) what is the function of scikit-network package?

- social network analysis
- analysis of large graphs

Q49) Write the name of the library from which the following methods originate:
i) find() — BeautifulSoup Date : _____
ii) get() — requests Page No. _____
iii) find-all() — BeautifulSoup
iv) prettyify() — BeautifulSoup

Q50) we are analyzing how much the no of lecture hours attended by students affect their exam scores. which excel function would you use as a starting point in this analysis?

SLOPE()

Q51) How to carry out a regression analysis in Excel ?
Data Analysis Toolpak.

- What is the function of scikit-network.
 - social network analysis
 - analysis of large graphs.
- How to get distance in km b/w two locations using geopy module?
geopy.distance.distance(loc1, loc2).km
- QD Joss? Inner Join.